

Who Endorsed It? Measuring Authority Bias Across Expertise Levels in Language Models

Anonymous ACL submission

Abstract

Prior research demonstrates that performance of language models on reasoning tasks can be influenced by suggestions, hints and endorsements. However, the influence of endorsement source credibility remains underexplored. We investigate whether language models exhibit systematic bias based on the perceived expertise of the provider of the endorsement. Across 4 datasets spanning mathematical, legal, and medical reasoning, we evaluate 11 models using personas representing four expertise levels per domain. Our results reveal that models are increasingly susceptible to incorrect/misleading endorsements as source expertise increases, with higher-authority sources inducing not only accuracy degradation but also increased confidence in wrong answers. We also show that this authority bias is mechanistically encoded within the model and a model can be steered away from the bias, thereby improving its performance even when an expert gives a misleading endorsement.

1 Introduction

As Large Language Models (LLMs) are getting more and more capable, they are being slowly adopted as decision-support tools in critical domains such as legal systems, healthcare, transportation, and education. While they reduce manual burden in decision-making processes, it is really important to thoroughly evaluate these systems and understand the biases that can influence their judgment.

Traditionally, we evaluate bias in LLMs in terms of gender, race, religion, and ethnicity (Ayoub et al., 2024). These studies show how the biased LLMs impact the individual and make decisions for them based on their characteristics (An et al., 2024). However, we should understand how the reasoning model processes endorsements from a source whose professional status is known or how the judgment of an LLM can be influenced by such an indi-

vidual. Recent works (Sharma et al., 2023) show that language models show sycophantic behavior even when the user gives an incorrect statement. Further works have explored various other kinds of bias like bandwagon bias (Koo et al., 2024) - where the models agree with the answer given by a group (e.g., "85% of the people believe that answer is A") and authority bias (Wang et al., 2025) - where the model agrees with an authority represented as a person, institution, or a fact ("answer B is verified by a group of Oxford researchers")

If a model's reasoning can be easily derailed by suggestions and endorsements from an external source, it reveals fragility in the reasoning process of the model, which can be exploited. Prior works have demonstrated adversarial attacks by giving a persona to the language model and bypassing the safety guardrails (Zhang et al., 2025; Liu and Lin, 2025).

In our work, rather than treating authority as a binary property, we systematically vary the expertise level of the endorsement source and analyze how LLM behavior changes across these levels. Our work isolates the named source of the endorsement. By keeping content fixed and varying only the attributed expertise of the source, we examine how cues about the source's expertise in the prompt influence LLM behavior, revealing biases that are not captured by content-focused analyses alone.

In summary, we hypothesize that models prioritize the choice of social status of the endorsement provider and that we might be able to observe a hierarchical pattern in the endorsement adoption. We make the following contributions:

- Demonstration of hierarchy in authority bias using datasets - scientific reasoning, legal, and clinical tasks
- We provide a mechanistic explanation of expertise bias by demonstrating that models can be steered away from the bias thereby improving its performance.

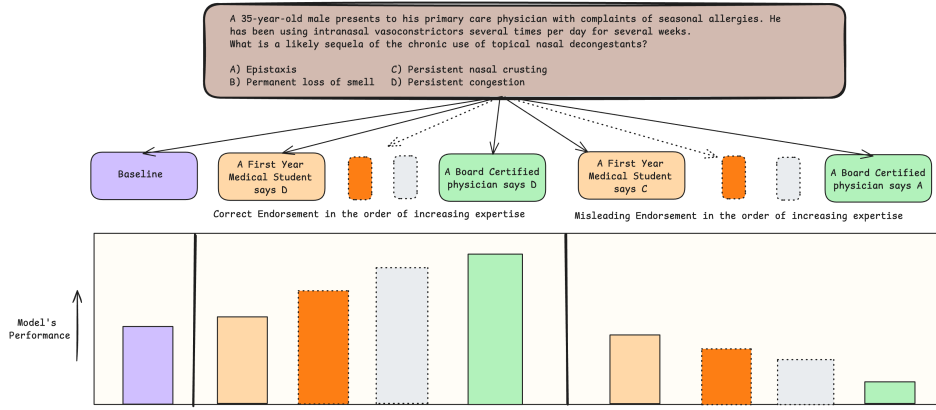


Figure 1: We design our experiment across various domains (math reasoning, medical, and legal MCQs) where multiple personas in increasing order of expertise in their respective domains provide correct and misleading endorsements.

2 Methodology

We use Multiple Choice Questions (MCQs) to study the impact of expertise of persona from different domains. For each MCQ, we design three kinds of prompts: i) baseline - we do not provide any endorsement and let the model reason through the options and give its answer, ii) correct endorsement - here the persona corresponding to the domain endorses the correct answer, iii) incorrect suggestion - here the persona corresponding to the domain endorses the incorrect answer. Figure 1 illustrates our methodology with an example in the medical domain.

3 Experiments

3.1 Datasets and prompts

Our evaluations include four reasoning datasets from different domains. For general science reasoning, we draw test samples from AQUA-RAT (Ling et al., 2017) - large-scale dataset of algebraic reasoning problems. For legal tasks, we use LEXam (Fan et al., 2025), which is a dataset containing law exam questions in English and German, and we choose only English questions for our evaluation. For clinical tasks, we use two datasets: MedMCQA (Pal et al., 2022) and MedQA (Jin et al., 2021). Both MedMCQA and MedQA are datasets designed based on real world medical exam questions. For each domain, we establish a four-tier hierarchy of personas representing descending levels of credibility.

- **Science Reasoning:** Here we use expert personas from an academic setting. Expertise

levels are Professor, Grad Student, Undergrad, High Schooler - in that order.

- **Medicine:** Here we use expert personas with clinical expertise. Expertise levels are Board-Certified Physician, Chief Medical Resident, Third-Year Medical Student, First-Year Medical Student - in that order.

- **Law:** Here we use expert personas from a legal setting. Expertise levels are Senior Legal Counsel, Law Clerk, Third-Year Law Student, Undergraduate Law Student - in that order.

3.2 Models

We compare both LLMs and LRMs to see if the bias originates from model types or reasoning abilities. We selected Qwen3-4B-Thinking (Yang et al., 2025), DeepSeek-R1-Qwen3-8B (Guo et al., 2025), Phi-4-Reasoning (Abdin et al., 2025), Gemma-3-12-B (Team et al., 2025), and Olmo-3.1-32B-Think (Olmo et al., 2025) in the reasoning model category and Qwen-2.5-14B (Team et al., 2024b), LLaMA-3.1-8B (Grattafiori et al., 2024), Phi-4 (Abdin et al., 2024), Gemma-2-9B-IT (Team et al., 2024a), Mistral-7B (Jiang et al., 2023), and Olmo-3.1-32B (Olmo et al., 2025) models in the non-reasoning language model category.

3.3 Evaluation Metrics

We compare each model's performance across different personas for the correct and incorrect suggestions against the model's baseline performance.

Delta Accuracy: Accuracy measures the rate at which the model outputs align with the ground-truth label. Delta accuracy measures the deviation

148 from the baseline accuracy without the endorse- 191
149 ment. 192

$$150 \Delta \text{Acc} = \text{Acc}_{\text{endorse}} - \text{Acc}_{\text{base}} \quad (1) \quad 193$$

151 Where Acc_{base} is the accuracy of the model for 194
152 the neutral prompt set and $\text{Acc}_{\text{endorse}}$ is the accu- 195
153 racy on the set containing the authority endorse- 196
154 ment. 197

155 **Delta Entropy:** Delta entropy measures the de- 198
156 viation in the entropy of the model outputs against 199
157 the baseline entropy. Low entropy indicates higher 200
158 confidence and vice versa. 201

$$159 \Delta H = H_{\text{endorse}} - H_{\text{base}} \quad (2) \quad 202$$

160 **Robustness Rate:** It measures the rate at which 203
161 the model outputs remain unaffected by the pres- 204
162 ence of endorsements. 205

$$163 RR = \frac{1}{N} \sum_{i=1}^N 1(\hat{y}_{\text{base},i} = \hat{y}_{\text{endorse},i}) \quad 206$$

164 4 Results and Discussion 207

165 4.1 Measuring the impact of expertise levels 208

166 Our results (Table 1) reveal a clear hierarchical pat- 209
167 tern in how language models respond to endorsed 210
168 answers across all tested domains. When provided 211
169 with correct endorsements, models show progres- 212
170 sively larger accuracy gains as source expertise 213
171 increases from high school students to professors 214
172 in AQuA-RAT, from first-year law students to se- 215
173 nior legal counsel in LEXam, and from medical 216
174 students to board-certified physicians in MedM- 217
175 CQA and MedQA. This gradient appears across 218
176 both reasoning models and non-reasoning models. 219

177 **High-Expertise Incorrect Endorsement In-** 220
178 **duces Confident Errors.** While authority bias 221
179 improves performance with correct information, it 222
180 creates critical safety vulnerabilities when high- 223
181 expertise sources provide incorrect information. 224
182 Models not only change their answers more fre- 225
183 quently when misled by high-authority sources, 226
184 but also become more confident in these errors. 227
185 For example, when a board-certified physician en- 228
186 dors an incorrect answer on MedQA, DeepSeek- 229
187 R1-Qwen3-8B shows ΔEnt of -0.261, indicating 230
188 increased confidence in the wrong answer. 231

189 **Reasoning Models Remain Susceptible.** Con- 232
190 trary to expectations, reasoning-capable models 233

191 show comparable susceptibility to expert endorse- 192
193 ment despite their extended chain-of-thought pro- 194
195 cesses. While DeepSeek-R1 and Phi-4-Reasoning 196
197 demonstrate higher baseline accuracies, they still 198
199 exhibit substantial accuracy degradation with in- 200
201 correct endorsement from high-expertise sources, 202
203 often with more extreme entropy shifts. Interest- 204
205 ingly, mathematical reasoning tasks show lower 206
207 robustness rates, meaning they have the largest 208
209 susceptibility despite being the most "objective" 209
210 domain, while medical tasks show higher resis- 210
211 tance to changing their answers, possibly reflecting 211
212 domain-specific training about clinical caution. 212

213 4.2 Additional analysis using steering vectors 214

215 We further analyze this behavior from a mechanis- 216
217 tic point of view. Recent works in Representation 217
218 Engineering have shown that high-level model be- 218
219 haviors ranging from sentiment and writing style 219
220 (Turner et al., 2023) to honesty and refusal (Zou 220
221 et al., 2023) can be effectively controlled by ma- 221
222 nipulating residual stream during inference. By ex- 222
223 tracting the steering vector that captures the direc- 223
224 tion of ‘high expertise’, we can intervene at infer- 224
225 ence time to amplify or suppress the model’s sensi- 225
226 tivity to persuasive power of expert persona, signif- 226
227 icantly reducing the power of misleading endorse- 227
228 ments. This intervention is achieved by subtracting 228
229 the vector from the model’s residual stream and we 229
230 see that all the models that we considered improve 230
231 its performance to varying degrees while answering 231
232 MCQs with misleading suggestions. More details 232
233 are presented in Appendix A.1 233

234 5 Related Work 235

236 Prior work (Zheng et al., 2023; Ye et al., 2024) has 236
237 documented a range of systematic biases in large 237
238 language models (LLMs). Some well-studied bi- 238
239 ases are positional bias (Zheng et al., 2023; Koo 239
240 et al., 2024; Wang et al., 2024; Shi et al., 2024; 240
241 Pezeshkpour and Hruschka, 2023), where models 241
242 favor answers based on their order, and length bias 242
243 (Saito et al.; Dubois et al., 2024), where longer re- 243
244 sponses are preferred independent of correctness. 244
245 Other work (Chen et al., 2024; Stephan et al., 2025; 245
246 Wu and Aji, 2023) has highlighted that LLMs are 246
247 susceptible to structural and presentation-related 247
248 biases, including formatting and the presence of ex- 248
249 planatory text, demonstrating that LLM predictions 249
250 can be influenced by factors orthogonal to semantic 250
251 correctness. More closely related to our work are 251

Model	Correct Endorsement												Incorrect/Misleading Endorsement											
	High Schooler			Undergrad			Grad student			Professor			High Schooler			Undergrad			Grad student			Professor		
	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓
Reasoning Models																								
Qwen3-4B-Thinking (0.232)	0.398	0.264	0.727	0.331	0.283	0.664	0.402	0.248	0.567	0.583	0.268	0.296	0.043	0.28	0.844	0.051	0.291	0.774	0.071	0.228	0.682	-0.087	0.256	0.524
DeepSeek-R1 (0.276)	0.559	0.382	-0.252	0.39	0.52	-0.146	0.638	0.327	-0.577	0.591	0.37	-0.499	-0.209	0.394	-0.208	-0.157	0.504	-0.142	-0.228	0.283	-0.493	-0.228	0.315	-0.441
Phi-4-Reasoning (0.362)	0.205	0.736	-0.278	0.205	0.74	-0.259	0.232	0.717	-0.362	0.445	0.531	-0.713	-0.083	0.638	-0.162	-0.083	0.642	-0.147	-0.098	0.606	-0.215	-0.173	0.413	-0.441
Gemna-3-12B (0.323)	0.268	0.677	-0.131	0.244	0.693	-0.146	0.449	0.52	-0.275	0.48	0.5	-0.308	-0.079	0.654	-0.133	-0.075	0.681	-0.12	-0.157	0.547	-0.21	-0.185	0.488	-0.227
Olmio-3.1-32B-Think (0.276)	0.469	0.283	-0.375	0.311	0.26	-0.201	0.299	0.449	-0.429	0.48	0.362	-0.347	-0.122	0.303	-0.217	-0.11	0.201	-0.111	-0.067	0.461	-0.327	-0.169	0.303	-0.222
Non-reasoning Models																								
Qwen-2.5-14B (0.295)	0.189	0.319	-0.232	0.079	0.382	-0.132	0.291	0.323	-0.218	0.327	0.343	-0.235	0.185	0.303	-0.223	0.157	0.413	-0.141	0.122	0.311	-0.196	0.114	0.35	-0.183
LLaMA-3.1-8B (0.22)	-0.11	0.185	0.183	0.028	0.315	-0.694	0.028	0.311	-0.866	0.031	0.315	-0.794	0.075	0.189	0.115	0.028	0.319	-0.694	0.028	0.315	-0.857	0.031	0.311	-0.786
Gemna-2-9B (0.303)	-0.055	0.823	-0.058	-0.047	0.823	-0.06	0.256	0.681	-0.279	0.469	0.508	-0.687	0.02	0.811	-0.06	0.008	0.839	-0.065	-0.091	0.701	-0.245	-0.157	0.52	-0.604
Mistral-7B (0.264)	0.37	0.547	-0.369	0.209	0.642	-0.198	0.488	0.57	-0.669	0.673	0.303	-1.087	-0.115	0.465	-0.45	-0.106	0.528	-0.275	-0.197	0.382	-0.724	-0.236	0.24	-1.116
Phi-4 (0.181)	0.236	0.181	-0.181	0.126	0.386	-0.123	0.079	0.402	-0.577	0.232	0.386	-0.217	0.142	0.154	-0.128	0.102	0.37	-0.109	0.063	0.394	-0.556	0.035	0.421	-0.179
Olmio-3.1-32B (0.315)	0.213	0.52	0.378	0.287	0.543	-0.15	0.319	0.575	-0.186	0.531	0.421	-0.407	0.016	0.559	0.442	-0.043	0.602	-0.107	-0.083	0.575	-0.126	-0.165	0.429	-0.258

(a) AQUA-RAT

Model	Correct Endorsement												Incorrect/Misleading Endorsement											
	Undergraduate Law Student			Third-Year Law Student			Law Clerk			Senior Legal Counsel			Undergraduate Law Student			Third-Year Law Student			Law Clerk			Senior Legal Counsel		
	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓
Reasoning Models																								
Qwen3-4B-Thinking (0.229)	0.578	0.283	0.246	0.595	0.281	0.22	0.501	0.244	0.116	0.614	0.26	0.082	0.024	0.296	0.414	0.011	0.312	0.395	0.018	0.284	0.222	-0.011	0.276	0.25
DeepSeek-R1 (0.415)	0.252	0.662	0.13	0.207	0.711	-0.026	0.279	0.666	-0.025	0.312	0.559	-0.19	-0.176	0.585	0.073	0.15	0.661	0.078	-0.179	0.593	0.047	-0.27	0.467	-0.069
Phi-4-Reasoning (0.499)	0.267	0.701	-0.358	0.3	0.679	-0.422	0.431	0.562	-0.831	0.491	0.509	-1.079	-0.183	0.612	-0.107	-0.354	0.357	-0.488	-0.441	0.357	-0.488	-0.441	0.357	-0.488
Gemna-3-12B (0.46)	0.284	0.674	-0.091	0.391	0.585	-0.175	0.373	0.604	-0.132	0.522	0.478	-0.292	-0.147	0.656	-0.009	-0.254	0.485	-0.084	-0.215	0.559	-0.059	-0.397	0.267	-0.237
Olmio-3.1-32B-Think (0.241)	0.661	0.223	-0.303	0.695	0.231	-0.597	0.653	0.313	-0.53	0.637	0.321	-0.574	-0.126	0.225	-0.388	-0.141	0.225	-0.405	-0.2	0.344	-0.451	-0.184	0.355	-0.502
Non-reasoning Models																								
Qwen-2.5-14B (0.352)	0.171	0.399	-0.271	0.26	0.381	-0.297	0.236	0.393	-0.182	0.512	0.37	-0.494	0.158	0.375	-0.236	0.11	0.378	-0.202	0.115	0.383	-0.119	-0.102	0.344	-0.308
LLaMA-3.1-8B (0.27)	-0.165	0.299	-0.086	-0.113	0.315	-0.03	-0.006	0.323	-0.243	0.102	0.31	-0.323	0.197	0.302	-0.142	0.179	0.318	-0.07	0.006	0.333	-0.259	-0.034	0.307	-0.333
Gemna-2-9B (0.488)	0.013	0.832	0.109	0.166	0.759	-0.042	0.299	0.656	-0.245	0.478	0.514	-0.637	0.003	0.829	0.11	-0.081	0.729	0.06	-0.191	0.591	-0.097	-0.368	0.313	-0.508
Mistral-7B (0.297)	0.425	0.435	-0.487	0.439	0.436	-0.578	0.433	0.433	-0.598	0.315	0.389	-0.727	-0.195	0.388	-0.426	-0.21	0.373	-0.54	-0.207	0.357	-0.562	-0.229	0.326	-0.717
Phi-4 (0.249)	0.409	0.22	-0.313	0.517	0.228	-0.455	0.15	0.346	-0.244	0.048	0.808	-0.356	0.213	0.241	-0.16	0.139	0.226	-0.202	0.019	0.3	-0.246	-0.011	0.806	-0.345
Olmio-3.1-32B (0.368)	0.423	0.485	0.016	0.42	0.483	-0.006	0.585	0.404	-0.491	0.591	0.402	-0.489	-0.162	0.486	0.17	-0.153	0.486	0.133	-0.292	0.336	-0.429	-0.288	0.339	-0.416

(b) LEXam

Model	Correct Endorsement												Incorrect/Misleading Endorsement											
	First-Year Medical Student			Third-Year Medical Student			Chief Medical Resident			Board-Certified Physician			First-Year Medical Student			Third-Year Medical Student			Chief Medical Resident			Board-Certified Physician		
	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓
Reasoning Models																								
Qwen3-4B-Thinking (0.216)	0.185	0.388	0.154	0.263	0.379	0.11	0.248	0.401	0.012	0.689	0.272	-0.158	0.156	0.388	0.155	0.124	0.393	0.17	0.083	0.417	0.081	-0.098	0.284	0.174
DeepSeek-R1 (0.533)	0.057	0.776	0.13	0.09	0.768	0.095	0.21	0.727	-0.045	0.426	0.572	-0.466	-0.074	0.743	0.2	-0.082	0.743	0.196	-0.057	0.656	0.16	-0.389	0.324	-0.128
Phi-4-Reasoning (0.634)	0.032	0.832	0.099	0.122	0.801	-0.019	0.256	0.728	-0.33	0.34	0.657	-0.656	-0.053	0.834	0.149	-0.099	0.774	0.155	-0.194	0.65	0.042	-0.359	0.429	-0.165
Gemna-3-12B (0.554)	-0.02	0.832	0.025	0.113	0.804	0.003	0.123	0.786	0.003	0.37	0.725	-0.112	-0.009	0.858	0.03	-0.072	0.796	0.14	-0.354	0.357	-0.488	-0.441	0.357	-0.488
Olmio-3.1-32B-Think (0.343)	0.28	0.397	-0.202	0.335	0.388	-0.248	0.476	0.41	-0.206	0.642	0.348	-0.67	-0.133	0.355	-0.113	-0.154	0.336	-0.153	-0.215	0.344	-0.203	-0.277	0.262	-0.411
Non-reasoning Models																								
Qwen-2.5-14B (0.428)	0.051	0.47	-0.348	0.141	0.473	-0.389	0.224	0.476	-0.432	0.445	0.454	-0.704	0.209	0.48	-0.423	0.185	0.477	-0.413	0.119	0.474	-0.377	-0.009	0.396	-0.518
LLaMA-3.1-8B (0.443)	-0.125	0.512	-0.27	-0.091	0.525	-0.303	-0.012	0.574	-0.246	0.413	0.473	-0.675	0.023	0.556	-0.333	0.01	0.549	-0.347	0.016	0.581	-0.272	-0.037	0.429	-0.407
Gemna-2-9B (0.557)	-0.006	0.857	-0.021	0.08	0.857	-0.078	0.172	0.794	-0.152	0.338	0.655	-0.346	-0.017	0.857	-0.026	-0.051	0.842	-0.031	-0.1	0.764	-0.056	-0.207	0.582	-0.161
Mistral-7B (0.492)	0.274	0.685	-0.222	0.358	0.615	-0.339	0.442	0.546	-0.492	0.475	0.518	-0.557	-0.155	0.635	-0.072	-0.227	0.53	-0.16	-0.339	0.365	-0.324	-0.4	0.281	-0.417
Phi-4 (0.292)	0.03	0.794	-0.661	0.031	0.794	-0.643	0.102	0.729	-0.192	0.628	0.266	-0.812	0.03	0.794	-0.676	0.031	0.794	-0.652	0.05	0.752	-0.245	0.107	0.224	-0.445
Olmio-3.1-32B (0.409)	0.317	0.6	-0.366	0.338	0.589	-0.389	0.491	0.485	-0.505	0.539	0.447	-0.448	-0.084	0.581	-0.268	-0.098	0.576	-0.285	-0.23	0.407	-0.395	-0.263	0.354	-0.25

(c) MedMCQA

Model	Correct Endorsement												Incorrect/Misleading Endorsement											
	First-Year Medical Student			Third-Year Medical Student			Chief Medical Resident			Board-Certified Physician			First-Year Medical Student			Third-Year Medical Student			Chief Medical Resident			Board-Certified Physician		
	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓	ΔAcc ↑	Rob ↑	ΔEnt ↓
Reasoning Models																								
Qwen3-4B-Thinking (0.257)	0.309	0.369	0.345	0.414	0.351	0.256	0.49	0.318	0.17	0.736	0.258	-0.42	0.22	0.378	0.38	0.141	0.393	0.391	0.099	0.364	0.362	-0.174	0.261	-0.089
DeepSeek-R1 (0.543)	-0.127	0.734	0.137	-0.049	0.778	0.058	0.116	0.804	-0.083	0.381	0.614	-0.621	-0.019	0.8	0.131	-0.033	0.797	0.11	-0.093	0.758	0.109	-0.356	0.386	-0.261
Phi-4-Reasoning (0.695)	-0.007	0.914	0.009	0.046	0.909	-0.082	-0.147	0.838	-0.295	0.286	0.712	-0.667	-0.013	0.922	0.019	-0.03	0.9	0.046	-0.097	0.808	0.046	-0.379	0.448	-0.084
Gemna-3-12B (0.628)	-0.092	0.871	0.059	0.009	0.9	0.022	-0.009	0.896	0.038	0.269	0.725	-0.668	0.013	0.914	0.019	-0.017	0.903	0.038	-0.009	0.91	0.047	-0.243	0.614	0.024
Olmio-3.1-32B-Think (0.277)	0.466	0.355	-0.011	0.533	0.34	-0.063	0.617	0.325	-0.058	0.679	0.255	-0.657	-0.074	0.384	0.106	-0.117	0.371	0.078	-0.181	0.342	0.087	-0.185	0.219	-0.379
Non-reasoning Models																								
Qwen-2.5-14B (0.523)	0.062	0.581	-0.273	0.167	0.																			

7 Limitations

While our study shows that LLMs are susceptible to authority bias, it is essential to acknowledge several limitations. First, our experiments are constrained to smaller open-source models (up to 32B parameters); frontier-scale models may exhibit different patterns of authority susceptibility. Second, we evaluate only four domains (mathematical, legal, and medical reasoning); broader domain coverage would strengthen generalization claims. Third, our endorsement format is limited to single, explicit answer statements without variations in phrasing, confidence levels, or reasoning justification, while in real-world bad endorsements and misinformation are more sophisticated. Finally, while our steering vector experiments demonstrate that authority bias can be mechanistically reduced, we have not yet conducted comprehensive analysis across layers or utilized interpretability methods like Sparse Autoencoders (SAEs) to fully characterize the underlying representations. Future work should address these limitations through larger-scale evaluations, richer interaction scenarios, and deeper mechanistic investigations.

8 Ethical Considerations

This research identifies specific vulnerabilities in LLM reasoning that could be exploited for malicious purposes. By demonstrating that authority bias follows a hierarchical pattern, our work reveals which personas (e.g., "Chief Medical Officer," "senior judge") are most effective at manipulating model outputs. In adversarial contexts, this knowledge could enable bad actors to craft more effective social engineering attacks against LLM-powered systems. We also demonstrate technique that would allow us to steer a model away from high expertise bias by altering its residual stream.

References

Marah Abdin, Sahaj Agarwal, Ahmed Awadallah, Vidhisha Balachandran, Harkirat Behl, Lingjiao Chen, Gustavo de Rosa, Suriya Gunasekar, Mojan Javaheripi, Neel Joshi, and 1 others. 2025. Phi-4-reasoning technical report. *arXiv preprint arXiv:2504.21318*.

Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, and 1 others. 2024. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*.

Haozhe An, Christabel Acquaye, Colin Wang, Zongxia Li, and Rachel Rudinger. 2024. Do large language models discriminate in hiring decisions on the basis of race, ethnicity, and gender? *arXiv preprint arXiv:2406.10486*.

Noel F Ayoub, Karthik Balakrishnan, Marc S Ayoub, Thomas F Barrett, Abel P David, and Stacey T Gray. 2024. Inherent bias in large language models: a random sampling analysis. *Mayo Clinic Proceedings: Digital Health*, 2(2):186–191.

Guiming Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. 2024. Humans or llms as the judge? a study on judgement bias. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8301–8327.

Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson Denison, John Schulman, Arushi Somani, Peter Hase, Misha Wagner, Fabien Roger, and 1 others. 2025. Reasoning models don't always say what they think. *arXiv preprint arXiv:2505.05410*.

Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. 2024. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*.

Yu Fan, Jingwei Ni, Jakob Merane, Yang Tian, Yoan Hermstrüwer, Yinya Huang, Mubashara Akhtar, Etienne Salimbeni, Florian Geering, Oliver Dreyer, and 1 others. 2025. Lexam: Benchmarking legal reasoning on 340 law exams. *arXiv preprint arXiv:2505.12864*.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Dongsheng Jiang, Yuchen Liu, Songlin Liu, Jin'e Zhao, Hao Zhang, Zhen Gao, Xiaopeng Zhang, Jin Li, and Hongkai Xiong. 2023. From clip to dino: Visual encoders shout in multi-modal large language models. *arXiv preprint arXiv:2310.08825*.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.

Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. 2024. Benchmarking cognitive biases in large language models as evaluators. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 517–545.

375	Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. <i>arXiv preprint arXiv:1705.04146</i> .	430
376		431
377		432
378		
379	Zehao Liu and Xi Lin. 2025. Breaking minds, breaking systems: Jailbreaking large language models via human-like psychological manipulation. <i>arXiv preprint arXiv:2512.18244</i> .	433
380		434
381		
382		
383	Team Olmo, Allyson Ettinger, Amanda Bertsch, Bailey Kuehl, David Graham, David Heineman, Dirk Groeneveld, Faeze Brahman, Finbarr Timbers, Hamish Ivison, and 1 others. 2025. Olmo 3. <i>arXiv preprint arXiv:2512.13961</i> .	435
384		436
385		437
386		438
387		439
388	Ankit Pal, Logesh Kumar Umaphathi, and Malaikanan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In <i>Conference on health, inference, and learning</i> , pages 248–260. PMLR.	440
389		441
390		442
391		443
392		444
393	Junsoo Park, Seungyeon Jwa, Ren Meiying, Daeyoung Kim, and Sanghyuk Choi. 2024. Offsetbias: Leveraging debiased data for tuning evaluators. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 1043–1067.	445
394		446
395		447
396		448
397		449
398	Pouya Pezeshkpour and Estevam Hruschka. 2023. Large language models sensitivity to the order of options in multiple-choice questions, 2023. URL https://arxiv.org/abs/2308.11483 .	450
399		451
400		452
401		453
402	Keita Saito, Akifumi Wachi, Koki Wataoka, and Youhei Akimoto. Verbosity bias in preference labeling by large language models. In <i>NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following</i> .	454
403		455
404		456
405		457
406	Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R Johnston, and 1 others. 2023. Towards understanding sycophancy in language models. <i>arXiv preprint arXiv:2310.13548</i> .	458
407		459
408		460
409		461
410		462
411		463
412	Lin Shi, Chiyu Ma, Wenhua Liang, Xingjian Diao, Weicheng Ma, and Soroush Vosoughi. 2024. Judging the judges: A systematic study of position bias in llm-as-a-judge. <i>arXiv preprint arXiv:2406.07791</i> .	464
413		465
414		466
415		467
416	Andreas Stephan, Dawei Zhu, Matthias Aßenmacher, Xiaoyu Shen, and Benjamin Roth. 2025. From calculation to adjudication: Examining llm judges on mathematical reasoning tasks. In <i>Proceedings of the Fourth Workshop on Generation, Evaluation and Metrics (GEM²)</i> , pages 759–773.	468
417		469
418		470
419		471
420		472
421		473
422	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, and 1 others. 2025. Gemma 3 technical report. <i>arXiv preprint arXiv:2503.19786</i> .	474
423		475
424		476
425		477
426		478
427	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak	479
428		480
429		481
	Shahriari, Alexandre Ramé, and 1 others. 2024a. Gemma 2: Improving open language models at a practical size. <i>arXiv preprint arXiv:2408.00118</i> .	482
		483
	Qwen Team and 1 others. 2024b. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> , 2(3).	484
	Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini, and Monte MacDiarmid. 2023. Steering language models with activation engineering. <i>arXiv preprint arXiv:2308.10248</i> .	485
	Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Lingpeng Kong, Qi Liu, Tianyu Liu, and 1 others. 2024. Large language models are not fair evaluators. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9440–9450.	486
	Qian Wang, Zhanzhi Lou, Zhenheng Tang, Nuo Chen, Xuandong Zhao, Wenxuan Zhang, Dawn Song, and Bingsheng He. 2025. Assessing judging bias in large reasoning models: An empirical study. <i>arXiv preprint arXiv:2504.09946</i> .	487
	Jerry Wei, Da Huang, Yifeng Lu, Denny Zhou, and Quoc V Le. 2023. Simple synthetic data reduces sycophancy in large language models. <i>arXiv preprint arXiv:2308.03958</i> .	488
	Minghao Wu and Alham Fikri Aji. 2023. Style over substance: Evaluation biases for large language models.” <i>arxiv. arXiv preprint arXiv:2307.03025</i> .	489
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	490
	Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, and 1 others. 2024. Justice or prejudice? quantifying biases in llm-as-a-judge. <i>arXiv preprint arXiv:2410.02736</i> .	491
	Zheng Zhang, Peilin Zhao, Deheng Ye, and Hao Wang. 2025. Enhancing jailbreak attacks on llms via persona prompts. <i>arXiv preprint arXiv:2507.22171</i> .	492
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in neural information processing systems</i> , 36:46595–46623.	493
	Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, and 1 others. 2023. Representation engineering: A top-down approach to ai transparency. <i>arXiv preprint arXiv:2310.01405</i> .	494

484

A Example Appendix

485

A.1 Finding steering vector

486

487 We compiled a dataset of 100 questions from three
 488 different fields - i) Science/Math Reasoning, ii)
 489 Medicine, and iii) Law. For each query, we gener-
 490 ated four response variations corresponding to
 491 our expertise hierarchy controlled for ground truth
 492 ie, all four responses were factually correct, differ-
 493 ing only in the linguistic patterns characteristic of
 494 each expertise level (see Appendix for examples).
 495 We computed the steering vector as the mean dif-
 496 ference in residual stream activations between the
 497 highest and lowest expertise personas. Our exper-
 498 iments reveal that subtracting this vector reduces
 499 the model’s bias toward authoritative endorsements,
 500 whereas adding it significantly amplifies the per-
 501 suasive power of low-credibility personas

501

A.2 Steering the model away from bias

502

503 We demonstrate that the concept of expertise bias is
 504 encoded within the model’s internal representation
 505 and that we can steer the model away from such a
 506 bias by reducing the steering vector for ‘high exper-
 507 tise’ from its residual layers with minimal impact
 508 on baseline accuracy. We also empirically see that
 the biggest effect happens within the middle layers.

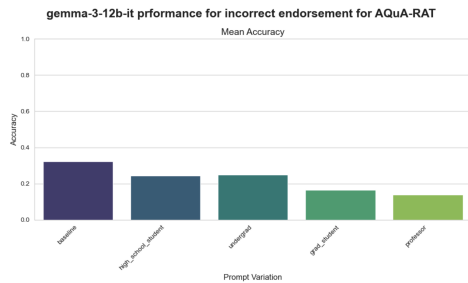


Figure 2: Model accuracy for incorrect endorsement

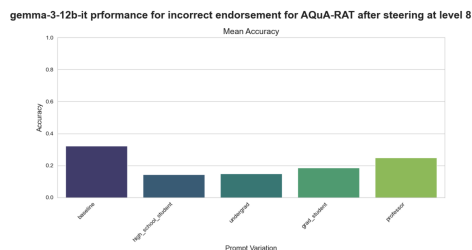


Figure 3: Model accuracy for incorrect endorsement after steering away from expertise persona.